



پروژه درس مقدمه‌ای بر یادگیری ماشین

نام و نام خانوادگی: امیرمهدی جعفری فشارکی

شماره دانشجویی: ۹۸۱۰۹۶۴۵

استاد درس: دکتر محمدزاده

تاریخ: ۱۴۰۰ بهمن ۴

فهرست مطالب

۱	توضیحات اولیه
۲	کلاس Person
۳	۱.۲ __init__ متد
۳	۲.۲ find_t_stimulated متد
۳	۳.۲ epoching متد
۴	۴.۲ create_all_features متد
۴	۵.۲ fisher_score_index متد
۴	۶.۲ cross_validation متد
۴	۷.۲ scorer متد
۵	۸.۲ find_best_model متد
۵	۹.۲ fit متد
۵	۱۰.۲ load_test_data متد
۵	۱۱.۲ find_test_word متد
۵	۱۱.۲ display_type1 متد
۵	۱۱.۲ display_type2 متد
۶	۳ نتایج به دست آمده برای هر شخص
۶	۱.۳ شخص ۱
۶	۲.۳ شخص ۲
۷	۳.۳ شخص ۳
۷	۴.۳ شخص ۴
۸	۵.۳ شخص ۵
۸	۶.۳ شخص ۶
۹	۷.۳ شخص ۷
۹	۸.۳ شخص ۸
۱۰	۹.۳ شخص ۹

۱ توضیحات اولیه

کد پروژه به همراه تمام داده ها و نتایج آن در فایل ML_Project_98109645.ipynb قرار دارند. این کد از چند بخش تشکیل شده است.

- در بلوک اول، تمامی کتابخانه و توابع لازم به کد اضافه شده اند. کتابخانه های استفاده شده در این پروژه عبارتند از: Pandas ، Numpy ، Scipy و Scikit learn ، Matplotlib
- در بلوک بعدی، داده های آموزش و تست در متغیر های datai_test و datai ذخیره شده اند که در حقیقت نشان دهنده شماره شخص می باشد.
- بلوک سوم کد که شامل پیاده سازی کلاسی به نام Person است در حقیقت مهمترین بخش پروژه می باشد. در واقع برای هر کدام از اشخاص، یک نمونه از این کلاس می سازیم و تمام اعمال مختلف نظری درست کردن ویژگی ها، انتخاب ویژگی ها با استفاده از Cross-Validation و Fisher Score ، انتخاب بهترین مدل، اندازه گیری دقت و پیش‌بینی کلمه برای داده تست و ... در این کلاس پیاده سازی شده است.
- در ادامه برای هر کدام از ۹ شخص، یک نمونه از کلاس ساخته شده و سپس بهترین مدل را برای آن پیدا و در نهایت کلمه متناظر با داده تست آن به دست آمده است.

در بخش های بعدی گزارش، ابتدا به پیاده سازی کلاس Person و سپس به نتایج به دست آمده برای هر کدام از افراد با جزئیات بیشتر پرداخته شده است.

۲ کلاس Person

این کلاس در حقیقت بخش اصلی پروژه را تشکیل داده و همان طور که پیشتر ذکر شد، برای هر شخص داده‌های آموزش به این کلاس داده شده و مدل مورد نظر برای آن شخص در درون این کلاس ذخیره شده. سپس به همین شکل با استفاده از یکی از توابع این کلاس، داده تست مورد نظر ورودی داده شده و برای آن، کلمات مورد نظر پیش‌بینی می‌شود. در ادامه توضیحات دقیق‌تر در رابطه با هر کدام از متد‌های مورد استفاده در این کلاس داده شده است.

۱.۲ متد __init__

این متد، در حقیقت فایل mat را که توسط تابع loadmat از کتابخانه io.scipy خوانده شده است را به عنوان ورودی می‌گیرد و سپس کل داده‌ها و همچنین زمان نمونه برداری را از روی آن استخراج می‌کند. این داده‌ها در متغیرهای data و t_sampling ذخیره می‌شوند.

۲.۲ متد find_t_stimulated

این متد، کل داده‌ها را به عنوان ورودی گرفته و آرایه‌ای از زمان‌هایی که در آن‌ها تحریک رخ داده است را به عنوان خروجی باز می‌گرداند. همچنین یک متغیر دیگر به عنوان test_data به عنوان ورودی به این تابع داده می‌شود که تعیین می‌کند که داده مورد نظر که قرار است زمان‌های تحریک روی آن به دست بباید داده تست است یا خیر. تفاوت این دو داده این است که تمام داده‌های ما، در حقیقت با هر بار تحریک، به اندازه ۴ نمونه تحریک مورد نظر روش می‌ماند و از طرفی ویژگی‌های این ۴ نمونه پشت سر هم به دلیل شباهت بسیار زیاد، عمل باعث می‌شود که به نوعی از هر داده در آموزش، ۴ کپی یکسان داشته باشیم. این موضوع در هنگام آموزش می‌تواند موجب به وجود آمدن overfit در داده‌ها و کاهش دقت در داده تست شود. به همین دلیل در داده آموزش یعنی زمانی که test_data با False برابر باشد، نمونه‌ها ۴ تا ۴ تا به خروجی داده می‌شود. به عبارت دیگر نمونه‌ها را با نرخ ۴ Downsample می‌کنیم.

۳.۲ متد epoching

این متد، داده‌ها و بردار زمان تحریک را به عنوان ورودی می‌گیرد و برای هر زمان تحریک و هر کanal، ۱۲۸ داده حول آن زمان از یک دهم قبل تا چهار دهم ثانیه بعد از زمان تحریک را ذخیره می‌کند و به عنوان یک آرایه با ابعاد تعداد نمونه در ۱۰۲۴ به عنوان خروجی باز می‌گرداند. همچنین باز هم یک متغیر test_data به عنوان ورودی داده می‌شود که صرفا در حالتی که داده آموزش باشد، علاوه بر این داده‌های زمانی، برچسب‌های مربوط به آن نیز به عنوان خروجی داده می‌شود.

٤.٢ متد `create_all_features`

این متد جزو مهمترین متدهای این کلاس است. در این متد، نمونه‌های زمانی استخراج شده توسط متد epoching به عنوان ورودی به این متد داده می‌شود و در خروجی، ویژگی‌های مربوط به آن نمونه‌های زمانی به عنوان خروجی داده می‌شوند. این ویژگی‌ها عبارتند از:

- میانگین گرفته شده ۱۲۸ داده روی ۸ کanal به عنوان یک داده ۱۲۸ تایی زمانی
- هیستوگرام برای هر ۸ کanal که توسط متد `create_hist_feature` ساخته می‌شود و داده‌ها را در بازه ۶۰-۱۲۰ به ۱۲ بخش تقسیم کرده و تعداد داده‌ای که در هر بخش قرار می‌گیرد را به عنوان یک ویژگی می‌دهد.
- میانگین داده‌های هر کanal
- واریانس داده‌های هر کanal
- همبستگی کanal‌ها با یکدیگر که توسط متد `create_correlation_feature` ۲۸ ویژگی برای هر نمونه تولید می‌کند.
- ویژگی‌های فرکانسی که توسط متد `create_frequency_features` تولید می‌شوند و عبارتند از:
 - انرژی ۵ بازه فرکانسی نیم تا ۴ هرتز، ۴ تا ۸ هرتز، ۸ تا ۱۳ هرتز، ۱۳ تا ۳۰ هرتز و ۳۰ تا ۴۰ هرتز به بعد
 - فرکانس میانگین
 - ۵ فرکانس با بیشترین پیک

٥.٢ متد `fisher_score_index`

این متد در حقیقت ماتریس متشكل از ویژگی‌ها و برچسب‌ها را به عنوان ورودی گرفته و برای هر ستون (که نشان دهنده هر ویژگی است) score را محاسبه کرده و به ترتیب از بیشترین به کمترین خروجی می‌دهد.

٦.٢ متد `cross_validation`

این متد، در حقیقت روی یک داده آموزش به ازای یک مدل مشخص و تعداد ویژگی مشخص (یعنی تعداد بهترین ویژگی‌ها بر اساس امتیاز فیشر)، cross validation انجام داده و امتیاز را برای مدل روی داده داده شده محاسبه و به عنوان خروجی می‌دهد.

٧.٢ متد `scorer`

این متد بر اساس متریکی که به عنوان ورودی به آن می‌دهیم، امتیاز را برای یک داده مشخص محاسبه می‌کند. همچنین لازم به ذکر است که متریک‌های قابل استفاده عبارتند از `accuracy` و `roc_auc`.

۸.۲ متد find_best_model

این متد، نمونه ها و ویژگی های مرتبط به آن را به همراه برچسب های مربوط گرفته. سپس ۲۰ درصد این داده ها را به عنوان تست و مابقی را به عنوان آموزش در نظر گرفته و از روی آن fisher score را برای هر ویژگی محاسبه و سپس روی مدل ها و تعداد ویژگی های متفاوت، cross validation انجام داده و بهترین مدل را (به همراه متغیر های دیگر نظیر اندیس های fisher score ، داده های تست و آموزش تقسیم شده و ...) به عنوان خروجی بازگردانده.

۹.۲ متد fit

این متد، ابتدا بهترین مدل را با استفاده از متد find_best_model که پیشتر معرفی شد یافته و سپس، امتیاز های متفاوت شامل accuracy و roc_auc و همچنین confusion matrix را برای بهترین مدل محاسبه و چاپ کرده. همچنین پس از آن، مدل را با تمام داده ها آموزش داده و این مدل را ذخیره کرده و برای این مدل نهایی نیز confusion matrix را محاسبه کرده.

۱۰.۲ متد load_test_data

این متد داده mat تست خوانده شده توسطتابع loadmat را به عنوان ورودی گرفته و ویژگی های آن، کارکتر های تحریک، زمان تحریک و دیگر پارامترهای مرتبط با آن را استخراج و ذخیره کرده.

۱۱.۲ متد find_test_word

این متد کلمه مورد نظر شخص در داده تست را پیشیبینی می کند. برای اینکار ابتدا تشخیص داده می شود که نحوه نمایش به شکل Single Character یا به شکل Row-Column می باشد. برای اینکار بیشترین عدد نمایش داده شده در ردیف ۱۰ ام چک شده و در صورتی که این عدد بیشتر از ۱۲ باشد، نحوه نمایش Single Character و در غیر این صورت به صورت Row-Column می باشد. پس از تشخیص نحوه نمایش، از یکی از دو متد display_type2 یا display_type1 استفاده شده تا کلمه مورد نظر بر اساس نحوه نمایش پیدا شود.

۱.۱۱.۲ متد display_type1

در صورتی که نحوه نمایش به صورت Single Character باشد، از این متد استفاده شده و این متد ابتدا کل زمان را به ۵ بخش تقسیم کرده و در هر بخش، حرفی که بیشترین تکرار داشته و برچسب مربوط به آن ۱ پیشیبینی شده است را به عنوان یک حرف در نظر گرفته و در نهایت با چسباندن این ۵ حرف به یکدیگر، کلمه نهایی را به دست می آورد.

۲.۱۱.۲ متد display_type2

این متد در زمانی که نحوه نمایش به صورت Row-Column باشد استفاده شده و نحوه کار آن این است که بین حالت هایی که تحریک دارای برچسب ۱ بوده یعنی واقعا شخص در حال نگاه کردن به آن کلمه بوده، حالت هایی که دقیقا پس از یک ستون، یک ردیف هم دارای برچسب ۱ بوده و یا بر عکس را در نظر گرفته و سپس بازه زمانی را به ۵ قسمت تقسیم و حرفی که بیشترین تکرار را داشته را به عنوان حرف مورد نظر در نظر گرفته و در نهایت با کنار هم قرار دادن این حروف، کلمه نهایی را ساخته و به عنوان خروجی باز می گردد.

۳ نتایج به دست آمده برای هر شخص

۱.۳ شخص ۱

- نحوه نمایش: تک حرفی

بهترین مدل: LDA با احتمال پیشین برابر

دقت: 80.93%

56.47% :roc_auc •

$$\begin{bmatrix} ۴۳۳ & ۹۶ \\ ۹ & ۴ \end{bmatrix}$$

ماتریس در هم ریختگی:

- کلمه نهایی: 0WBEW

۲.۳ شخص ۲

- نحوه نمایش: تک حرفی

بهترین مدل: لاجیستیک

دقت: 72.04%

47.47% :roc_auc •

$$\begin{bmatrix} ۳۸۵ & ۱۳۶ \\ ۱۵ & ۴ \end{bmatrix}$$

ماتریس در هم ریختگی:

- کلمه نهایی: 1UKAS

۳.۳ شخص ۳

- نحوه نمایش: سطري-ستوني

بهترین مدل: LDA با احتمال پيشين برابر

دقت: 69.44%

65.50% :roc_auc •

$$\begin{bmatrix} 106 & 42 \\ 13 & 19 \end{bmatrix}$$

- ماتريس در هم ریختگی:

کلمه نهايی: 7O2AD

۴.۳ شخص ۴

- نحوه نمایش: سطري-ستوني

بهترین مدل: جنگل تصادفي با ماکسيمم عمق ۴ و ۹ نخمينگر

دقت: 73.89%

55.18% :roc_auc •

$$\begin{bmatrix} 125 & 26 \\ 21 & 8 \end{bmatrix}$$

- ماتريس در هم ریختگی:

کلمه نهايی: 52ZCD

۵.۳ شخص ۵

- نحوه نمایش: سطري-ستوني

بهترین مدل: LDA با احتمال پيشين برابر

دقت: 63.33%

52.15% :roc_auc •

$$\begin{bmatrix} 104 & 49 \\ 17 & 10 \end{bmatrix}$$

ماتريس در هم ریختگی:

كلمه نهايى: XAPZT

۶.۳ شخص ۶

- نحوه نمایش: سطري-ستوني

بهترین مدل: LDA با احتمال پيشين برابر

دقت: 60.00%

50.25% :roc_auc •

$$\begin{bmatrix} 96 & 49 \\ 23 & 12 \end{bmatrix}$$

ماتريس در هم ریختگی:

كلمه نهايى: MJD4R

٧.٣ شخص ٧

• نحوه نمایش: سطري-ستوني

• بهترین مدل: جنگل تصادفي با ماکسیمم عمق ۶ و ۹ نخمنگر

• دقت: 70.56%

48.04% :roc_auc •

$$\begin{bmatrix} 124 & 16 \\ 37 & 3 \end{bmatrix}$$

• ماتریس در هم ریختگی:

• کلمه نهایی: PXLMU

٨.٣ شخص ٨

• نحوه نمایش: سطري-ستوني

• بهترین مدل: LDA با احتمال پیشین برابر

• دقت: 67.78%

58.36% :roc_auc •

$$\begin{bmatrix} 108 & 40 \\ 18 & 14 \end{bmatrix}$$

• ماتریس در هم ریختگی:

• کلمه نهایی: 3RV4D

۹.۳ شخص ۹

• نحوه نمایش: سطري-ستوني

• بهترین مدل: LDA با احتمال پيشين برابر

• دقت: 66.11%

52.96% :roc_auc •

$$\begin{bmatrix} 108 & 37 \\ 24 & 11 \end{bmatrix}$$

• ماتريس در هم ريختگی:

• کلمه نهايی: 5PZAN